

# ARC's SIMD Extensions for Multimedia Applications



Nigel Topham  
Chief Architect

# Quest for an Efficient High-Performance Embedded Solution

- Embedded media-rich applications need greater performance
- Cannot rely on increasing clock speed
- Power consumption is a critical constraint
- Parallelism is the future – but what form will it take?

~~VLIW~~

~~Superscalar~~

~~Simultaneous  
Multithreading  
(SMT)~~

☒ Data-level  
Parallelism

☒ Multiprocessor  
Parallelism

## Why Not?

- Compiler complexity
- Issue bandwidth
- Die-area bloat
- Poor sustained performance
- High power consumption

## Why Not?

- Hardware complexity
- Modest gains
- Not scalable
- May play a part, but is not itself a solution

## Why Not?

- Offers a throughput gain only
- May make sense for a large, inefficient core
- More effective to build multicore ARC systems

## Requirements

- ☒ SIMD, short vector
- ☒ Multiple instruction issue
- ☒ Configurable
- ☒ Application-specific
- ☒ High memory bandwidth

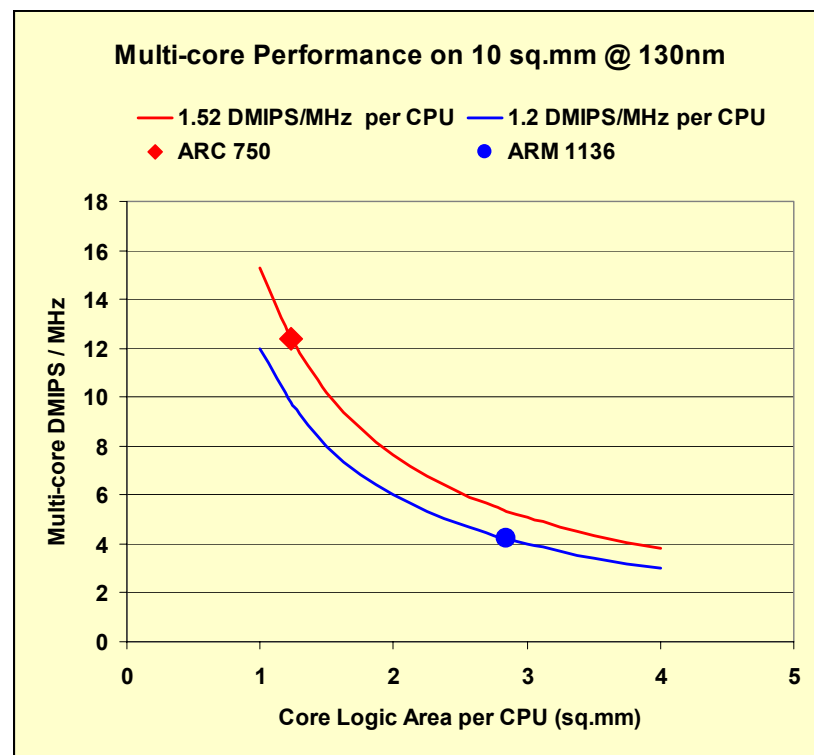
## Requirements

- ☒ SMP
- ☒ Heterogeneous
- ☒ Configurable
- ☒ Application-specific
- ☒ Inter-processor communications

- ARC's multicore history
  - 40% of ARC-based deployments involve multiple cores
  - Average number of cores per design = 6
- Small, fast, low-power processing elements are key
- Latest ARC™ 750D core is ideal for multicore implementations

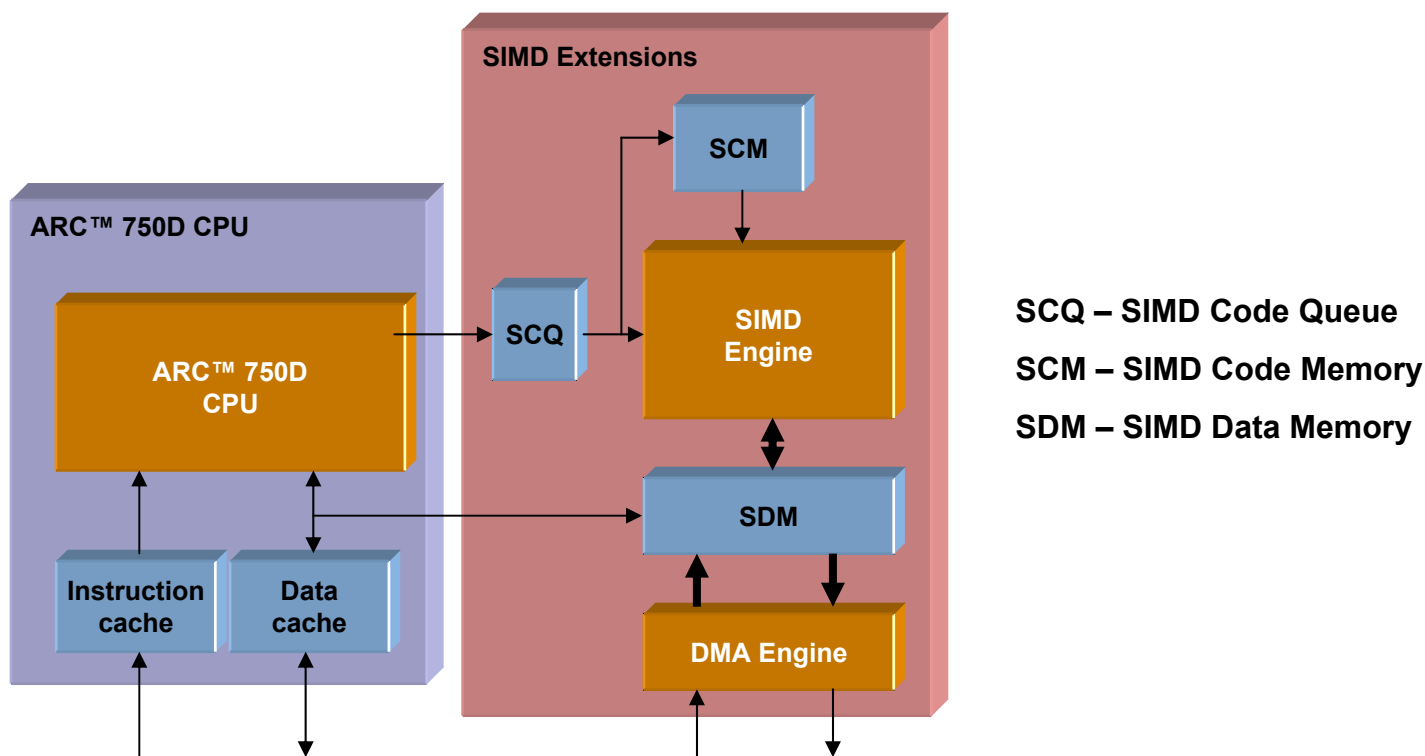
ARC™ 750D Performance		
	130 nm	90 nm
Frequency	533 MHz	700 MHz
DMIPS / MHz	1.526	1.526
DMIPS / CPU	813	1068

- ARC's multicore roadmap
  - Symmetric multiprocessor (SMP)
  - Heterogeneous scalar / SIMD



# Introducing ARC's SIMD Extensions

- High-performance SIMD extensions extend ARC™ 750D core
- 128-bit data-parallel operations for high-throughput processing
- Decoupled-control concept yields low-cost superscalar performance
- Media-oriented dual-channel two-dimensional DMA engine
- All delivered as ARCompact™ compatible ISA extensions



- Extends ARCompact™ ISA
- 128-bit vectors
  - 4 x 32 bits, 8 x 16 bits, 16 x 8 bits
  - Type depends on instruction
- Up to 64 vector registers
  - 24 in first implementation (vr00 – vr23)
- Scalar 16-bit registers (i0 – i7)
  - Eight scalar data or address values
  - Synonymous with vector register 0
- Vector accumulator
  - Third operand and result for some ops
  - 32-bit or extended 40-bit elements

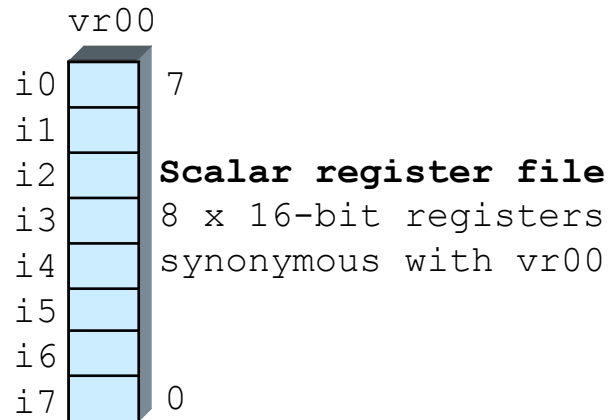
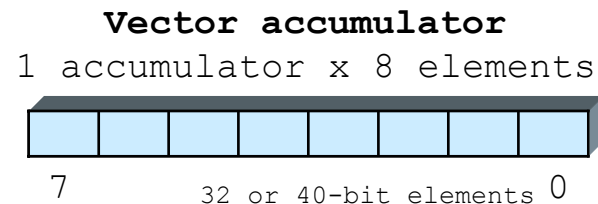
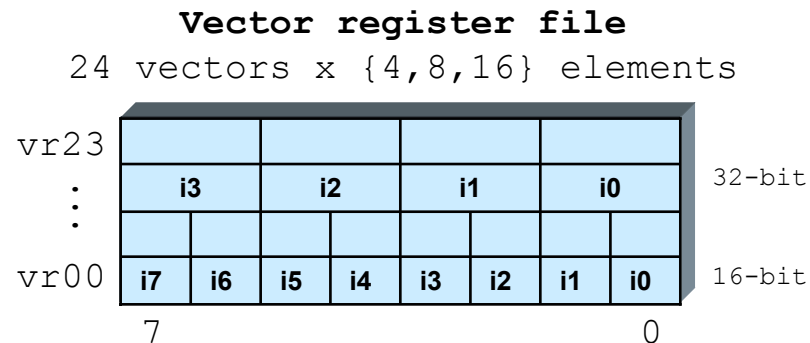
**Broadcast multiply with accumulation**

```
vbmulaw vr08, vr02, r3
```

**Semantics**

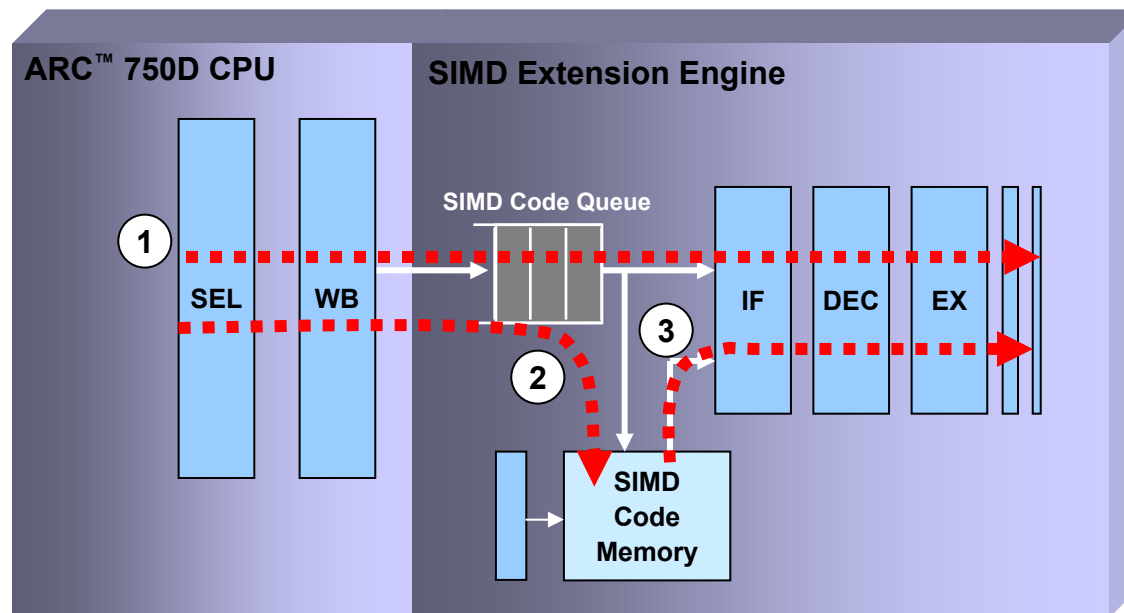
```
for (i = 0; i < 8; i++)
```

```
    vr08[i] = acc[i] = acc[i] + r3*vr02[i]
```



# Decoupling Scalar and SIMD Execution

- Closely coupled and decoupled execution modes



- ① **Closely coupled execution** — instructions issued directly from ARC™ 750D core
- ② **Macro record** — instructions fetched by 750D core and stored in SCM
- ③ **Decoupled execution:**
  - 750D core issues VRUN with macro start address in SCM
  - SIMD engine issues instructions autonomously from SCM
  - Thereafter, 750D core and SIMD engine operate in parallel



# Record / Replay Enables Parallel Execution

- Coupled model embeds SIMD code with ARC™ 750D core code
- VRUN initiates macro sequence
- Many VRUNS can be queued

## Tightly coupled execution model

```
example_loop:
    push_s    r13
    mov_s     r13,r6
    vmovzwb   vr00,r1,1
    vmovwb    vr00,r2,2
    vmovwb    vr00,r0,4
    vmovzwb   vr23,r4,1
    vmovwb    vr23,r5,2
    vmovwb    vr23,r3,4
    cmp_s     r13,0
    mov.gt    lp_count,r13
    lpgt      .L01
    vld128    vr01,[i0]
    vld128    vr02,[i0,0x10]
    vmrb      vr01,vr02,vr01
    vld128    vr03,[i1]
    vavrb     vr01,vr01,vr03
    vst128    vr01,[i2]
    vaddw     vr00,vr00,vr23
.L01      j      r13
```

## Decoupled execution model

```
example_loop:
    push_s    r13
    mov_s     r13,r6
    vmovzwb   vr00,r1,1
    vmovwb    vr00,r2,2
    vmovwb    vr00,r0,4
    vmovzwb   vr23,r4,1
    vmovwb    vr23,r5,2
    vmovwb    vr23,r3,4
    cmp_s     r13,0
    mov.gt    lp_count,r13
    lpgt      .L01
    vrun      r7
.L01      j      r13
```

## Macro sequence pre-recorded in SCM

```
[r7]:  vld128    vr01,[i0]
        vld128    vr02,[i0,0x10]
        vmrb      vr01,vr02,vr01
        vld128    vr03,[i1]
        vavrb     vr01,vr01,vr03
        vst128    vr01,[i2]
        vaddw     vr00,vr00,vr23
```

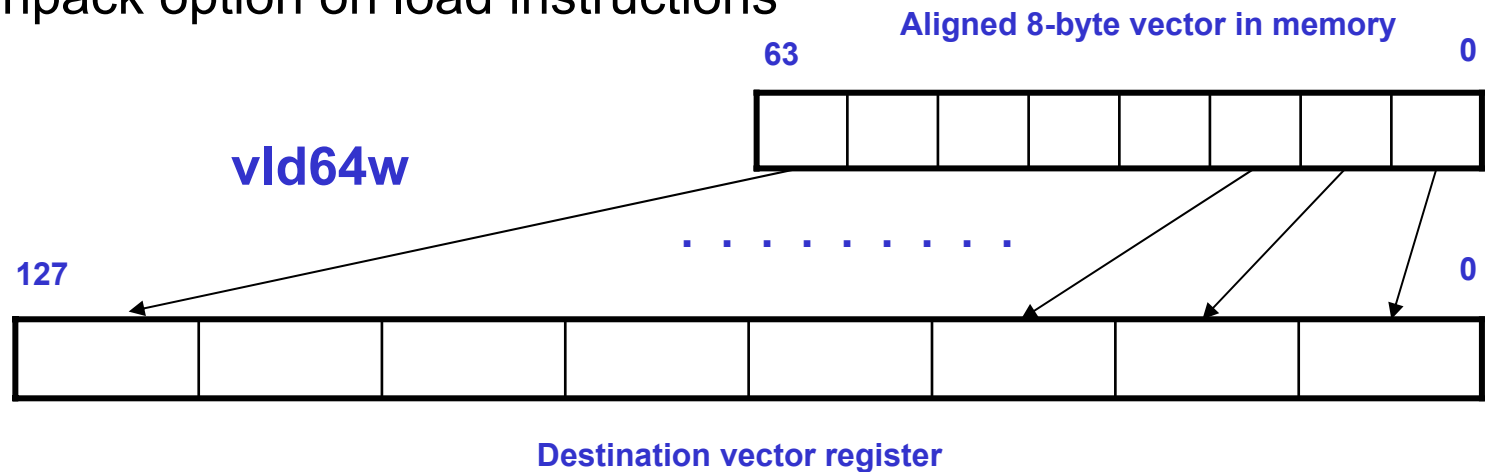
# SIMD Instruction-Set Overview

- 104 new vector instructions extend ARCompact™ ISA
- Carefully selected to optimize media application performance

Load	Store	Move	Logic			Arith			Compare
vld32	vst16	vmwv	vmvaw	vand	vandaw	vaddw	vbadw	vaddaw	veqw
vld64	vst32	vmvzw		vor		vsubw	vbsubw	vsubaw	vnew
vld64w	vst64	vmoww	vmovaw	vxor	vxoraw		vbrsubw		vlew
vld32wl	vst128	vmovzw		vbic	vbicaw	vsummw			vltw
vld32wh	vst128r					vaddsuw			
vld128						vavb			
vld128r						vavrb			
						vmulw	vbmulw	vmulaw	
						vmulfw	vbmulfw	vmulfaw	
								vbmulaw	
						vmaxw	vbmaw	vmaxaw	
						vminw	vbmaw	vminaw	
						vabsw		vabsaw	
						vdifw		vdifaw	
						vsignw			
Permute	Align	Scale	Pack	Unpack	Misc	Special	Macros	DMA	
vexch1	vmr1w	vmr1aw	vasrw	vasrpwb	vupbw	vnop	vd6tapf	vrec	vdirun
vexch2	vmr2w	vmr2aw	vasrrw	vasrrpwb	vupsbw	vint	vh264ft	vendrec	vdorun
vexch4	vmr3w	vmr3aw	vasrsrw		vupbaw		vh264f	vrun	vdiwr
	vmr4w	vmr4aw			vupsbaw		wc1ft	vrecrun	vdowr
	vmr5w	vmr5aw					wc1f		vdird
	vmr6w	vmr6aw							vdord
	vmr7w	vmr7aw							
	vsr8	vsr8aw							
	vmrb								

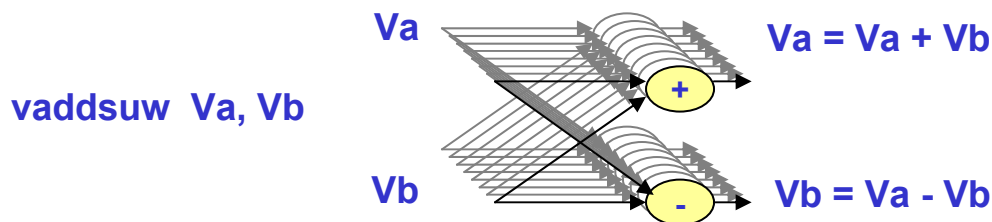


- Vector load or store up to 128 bits per cycle from SDM
- DMA moves data between SDM and external RAM
- Unpack option on load instructions



- Move vector elements (vector-to-vector move)
- Replicate scalar value (scalar-to-vector broadcast)
- Selective element assignment, determined by bit mask
- Deselected elements may be cleared or left intact

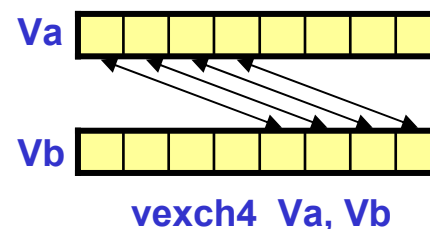
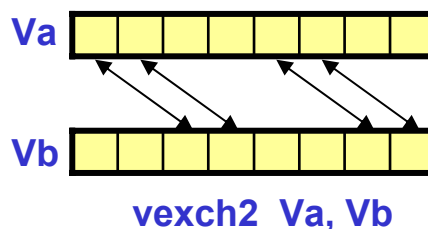
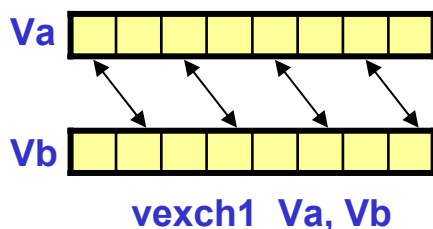
- Two-point vectorized butterfly – ideal for DCT and iDCT



- Vector multiplication
  - Vector x vector, or vector x scalar
  - Eight 16 x 16 multipliers
  - 16-bit integer or 1Q15 fractional results with rounding
- Byte-vector averaging for motion compensation
  - Half- and quarter-pixel motion vectors require block interpolation
  - Capable of 16 byte-average calculations per cycle
  - Hence,  $h$  cycles to interpolate a  $16 \times h$  block

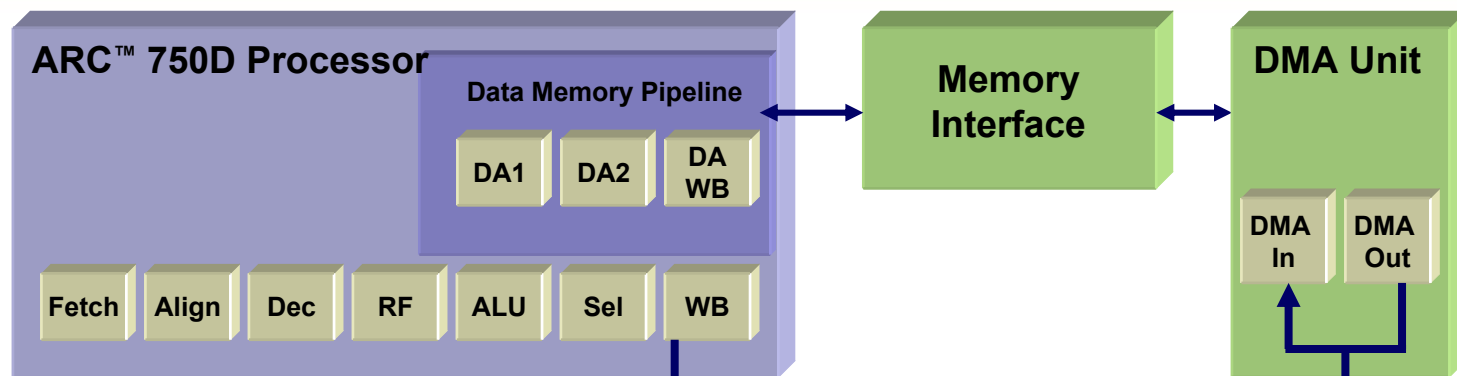
# Permutations and Matrix Transposition

- 2D block transforms have horizontal and vertical passes
- Matrix transposition occurs between passes (4 x 4 or 8 x 8)
- Permutation instructions yield fast transposition

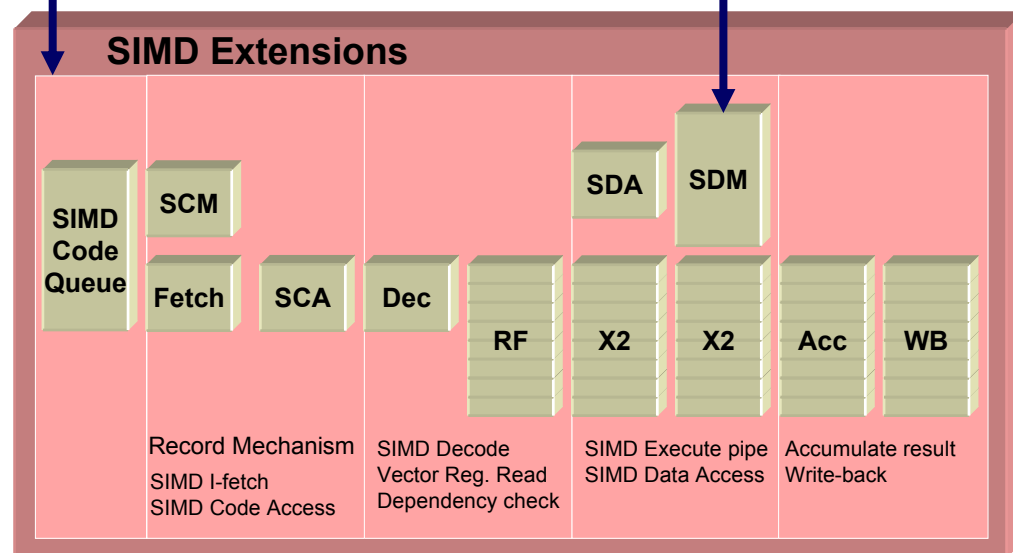


- Throughput:
  - One 4 x 4 matrix transposed in 4 instructions
  - Four 4 x 4 matrices transposed in 8 instructions
  - 8 x 8 matrix transposed in 12 instructions

# SIMD Extensions Pipeline



- SIMD instructions extend pipeline
- Operates at ARC™ 750D core clock (533 MHz)
  - Can run at lower speed if required and at different speeds from the ARC 750D
  - High parallelism makes this feasible
  - Reduced power consumption
- 128-bit datapaths
  - Optional 160-bit accumulator datapath
- 128-bit load / store path to SDM
- 8.5GByte/sec bandwidth DMA
- Decoupled control
  - Independent SIMD I-fetch from SCM
  - Executes SIMD basic blocks with VRUN
  - 750D core queues future blocks in SCQ
  - 750D core “gets ahead” of SIMD engine





# ARC™ Multimedia Subsystems Exploit SIMD Extensions

- Multistandard software-based platforms for audio/video media processing
- ARC™ Media Subsystem = ARC Video + ARC Sound Subsystem technologies
  - ARC Video: MPEG 1, 2, & 4, H.264, VC-1...and any future codecs
  - ARC Sound: All popular audio codecs
- Wide performance range
  - Low-frequency operation for mobile low-power applications
  - 600 core family sufficient for ARC Sound
  - SIMD extensions are powerful enough for H.264 and VC-1 decode at D1 resolution
- ARC Sound Advanced Subsystem
  - 12-channel AAC or MP3 encode, running on SIMD extensions

	Dedicated Hardware	General-Purpose CPU	ARC™ Media Subsystem
Flexibility	None	Unlimited	Unlimited
Die area	Additive per codec	Fixed but high	Fixed and low
Power consumption	Low	High, due to high clock frequency	Low
Operating frequency	Very low	Very high	Low

- The compelling case for software-based media processing
- Combining SIMD, multicore, and custom extensions is key

- DCT based on Arai, Agui, and Nakajima's scaled DCT<sup>1</sup>
  - Originally 8-point 1-D transform required 5 multiplies and 29 adds
  - ARC™ Multimedia subsystem implementations exploit vector butterfly extensively
  - Computes 8 x 8 1-D transform in 24 instructions
  - Full 2-D 8 x 8 DCT requires only 88 instructions
  - Quantization requires a further 8 instructions
  - 9.4x speedup due to SIMD extensions, compared with basic 750D core
- iDCT uses modified version of Loeffler's algorithm<sup>2</sup>
  - Originally 8-point 1-D transform required 11 multiplies and 29 adds
  - ARC Multimedia subsystems use 16 multiplies and 16 other ops (add, sub, butterfly)
  - Computes 8 x 8 1-D transform in 32 instructions
  - Full 2-D 8 x 8 IDCT requires only 116 instructions
  - 10.4x speedup due to SIMD extensions, compared with basic 750D core
  - Easily satisfies IEEE 1180 accuracy test

---

1. Arai, Agui, and Nakajima, Trans. IEICE E-71 (11): 1095

2. Loeffler, Ligtenberg, and Moschytz, Proc. ICASSP '89, pp. 988-991

- Typically, integer or fractional de-scaling is via arithmetic right shift
- However, accurate rounding is critical to numerical fidelity
- No rounding
  - Compute as:  
 $(x \gg d)$
  - Mean error =  $-0.5 + 0.5 * 2^{-d}$
- Simple rounding
  - Compute as:  
 $(x + (1 \ll (d-1))) \gg d$
  - Mean error =  $0.5 * 2^{-d}$
- Symmetric rounding in ARC™ SIMD extensions
  - Compute as:  
 $(x < 0) ? (x + (1 \ll (d-1)) - 1) \gg d : (x + (1 \ll (d-1))) \gg d$
  - Mean error is ZERO
  - Implemented in ARC Media subsystems right shift
  - Enables 16-bit fast DCT / iDCT to meet IEEE 1180 compliance

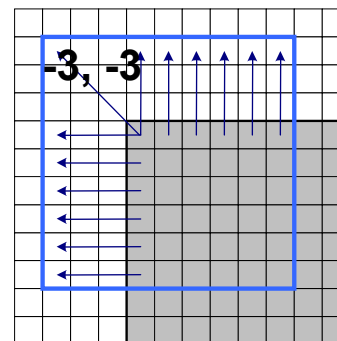


- SIMD ISA extended to support H.264 and VC-1
- Complex filtering operations split across two instructions
  1. Test deblocking conditions and set deblocking values
  2. Apply deblocking values to block boundary pixels
- H.264 and VC1 deblocking filters
  - `vh264ft, vvc1ft` – horizontal filter test
  - `vh264f vvc1f` – horizontal filter calculation
- Key optimization for H.264 and VC1 software implementation
  - Hides control complexity
  - Exposes irregular pixel-level parallelism to hardware
  - Speedup of 15x – 25x, exceeding the 8x expected of 8-way SIMD

# Custom Scalar Extensions: Entropy Decoding

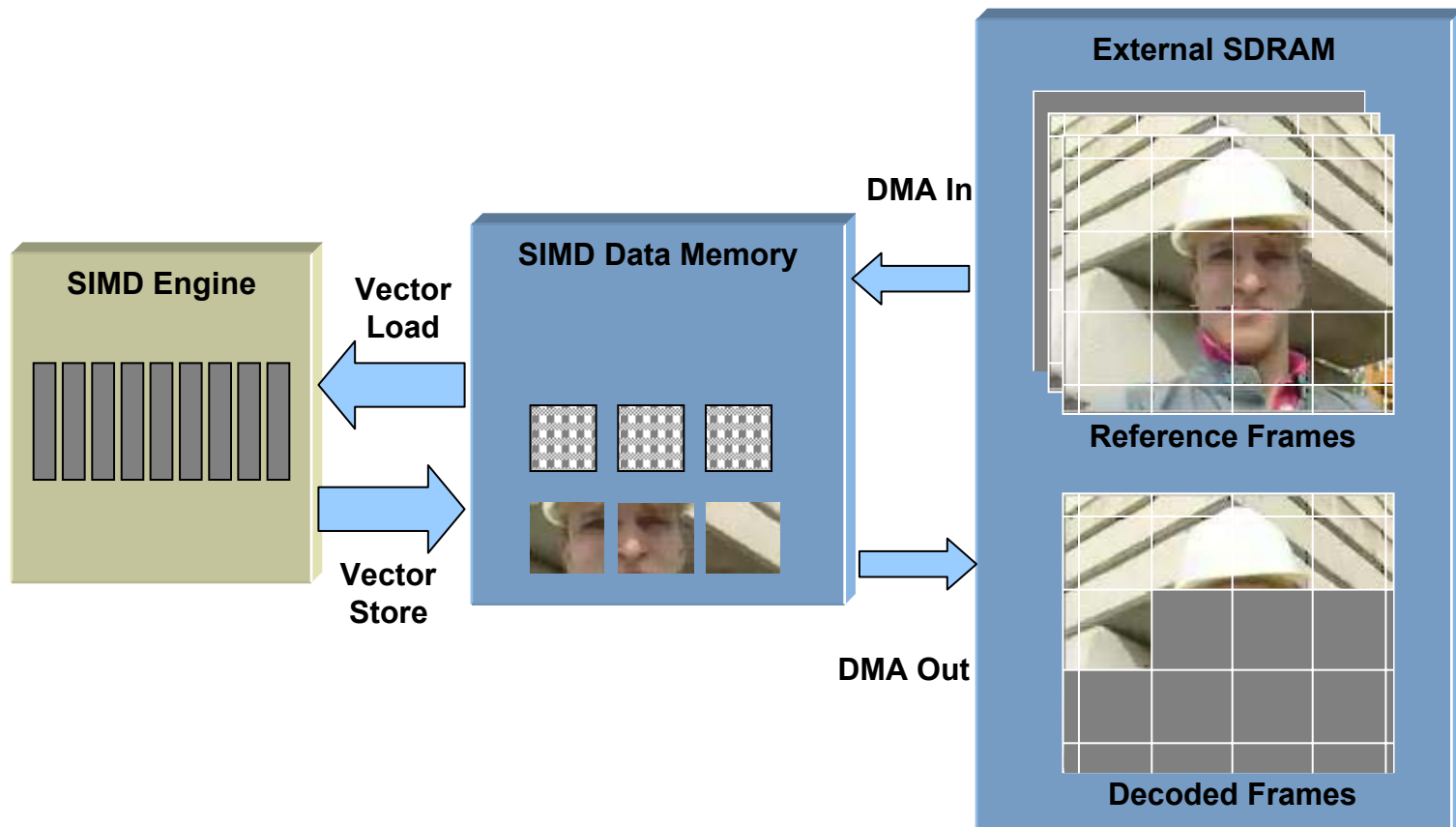
- Media codecs each define a bit-stream syntax
- Symbols encoded using a range of techniques
  - Huffman coding (JPEG, MPEG, VC-1)
  - Arithmetic coding (H.264)
  - Most complex schemes are context adaptive
- Decoding is a serial process – cannot be parallelized
- Goal is to decode one symbol per clock period (or instruction)
- ARC's scalar extension instructions provide ideal mechanism
- ARC solution:
  - Implement a single-cycle hardware decoder via extension instructions
  - All relevant codecs supported directly
    - VC-1, MPEG-1, 2, and 4 VLC and RVLC tables
    - H.264: CAVLC, CABAC, Exp-Golomb algorithms
- Major benefits:
  - No software VLC tables – no D-cache pollution from sparse structures
  - No D-cache dependency – predictable decode delays
  - Low die area cost
  - Extremely fast

- Lightweight command interface – fast DMA initiation
- Frame descriptors held in SDM
- DMA commands specify frame and (x,y) location in frame
- Clipping and filling
  - Blocks can extend outside frame
  - DMA fills or clips appropriately
  - A feature of H.264 and VC-1
- DMA is SDRAM aware
  - Address reordering to minimize page-boundary crossing
- Supports both interlaced and progressive scan

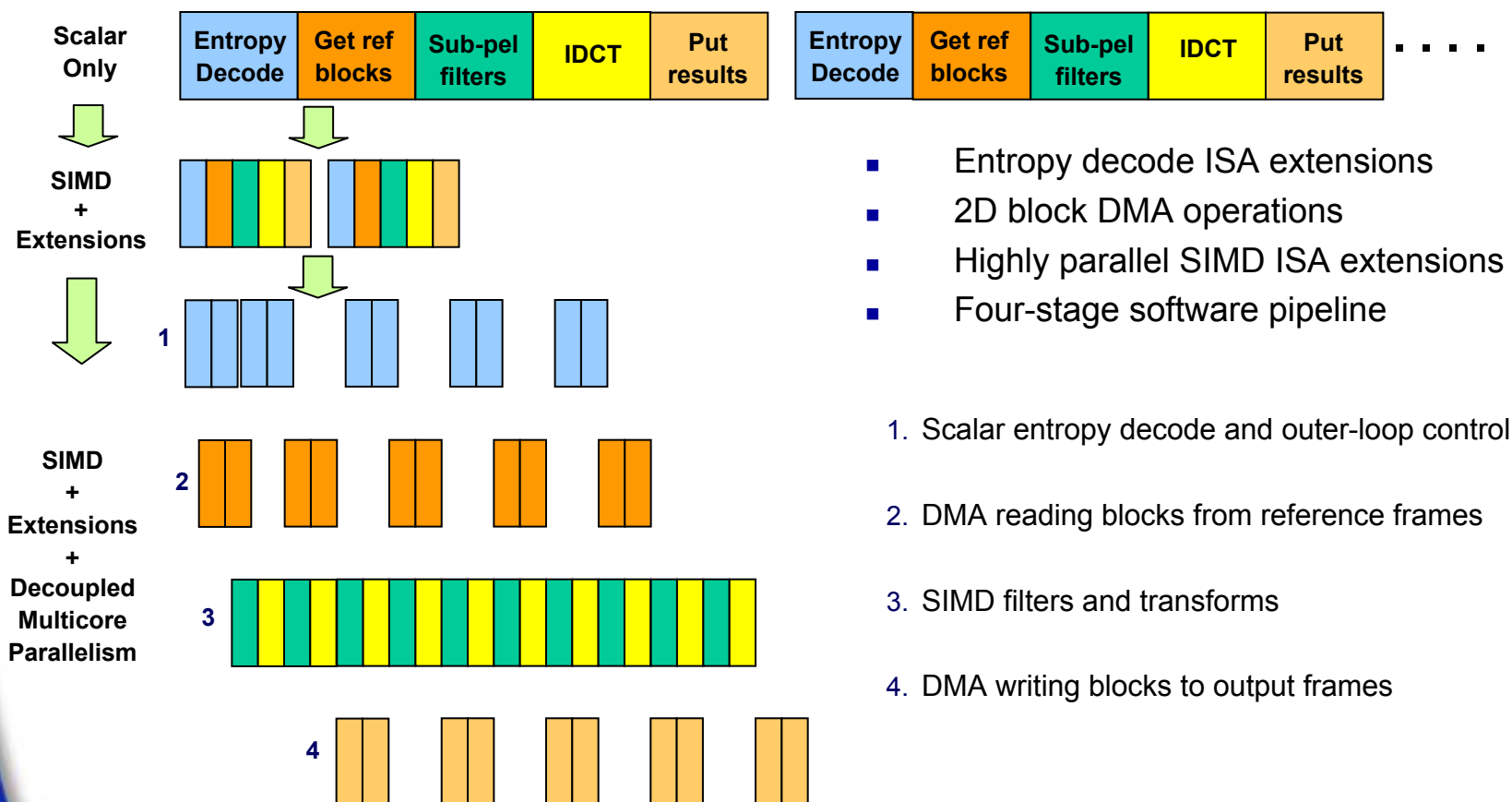


# Scalar, SIMD, and DMA = 3-Way Parallelism

- Vector loads and stores address SDM (fast, on-chip RAM)
- Reference pictures typically stored in external SDRAM
- DMA moves 2D blocks between SDM and external SDRAM



- Combines parallel processing with ISA extensions
  - Parallelism among three major units
  - ISA extensions for entropy decode and vector processing



# SIMD Video Decode Performance

- Exceptional performance for MPEG, H.264, and VC-1
- High sustained levels of data parallelism
- Illustrated with selection of key video computations
  - Coded and fully functional on ARCCangel™ 4 development platform
  - Cycle times include block load / store as well as SIMD computation

Codec	Computation	SIMD clocks	Block	Clocks / pixel-op	SIMD Speedup
MPEG-1 MPEG-2 MPEG-4	Forward DCT 8x8	113	8 x 8	1.77	<b>9.4</b>
	Inverse DCT 8x8	129	8 x 8	2.02	<b>10.4</b>
	Matrix transpose 8x8	12	8 x 8	0.19	<b>12.5</b>
	1/2-pel interpolation 16x16	83	16 x 16	1.30	<b>9.8</b>
MPEG-4	Vertical lowpass filter	453	16 x 16	7.08	<b>12.9</b>
	Horizontal lowpass filter	516	16 x 16	8.06	<b>9.8</b>
H.264	Integer transform	82	4x4	5.13	<b>10.3</b>
	<i>Inter-prediction</i>				
	vertical filtering	71	8x8	1.11	<b>25.2</b>
	horizontal filtering	116	8x8	1.81	<b>15.4</b>
	1/4-pel interpolation	46	8x8	0.72	<b>8.3</b>
	Intra-prediction (planar)	240	16x16	0.94	<b>7.7</b>
	<i>De-blocking filter</i>				
	horizontal filtering	2	6	0.33	<b>70.0</b>
	vertical filtering	54	24	2.25	<b>10.4</b>
VC-1	Integer transform	143	8 x 8	2.23	<b>10.7</b>
	Bilinear block filter	93	8 x 8	1.45	<b>11.0</b>
	Bicubic 2D filter	216	8 x 8	3.38	<b>11.0</b>



- SIMD extensions add 150K gates
- Shares ARC™ 750D core's low-power characteristics
- Capable of >500 MHz operation
  - Codecs typically require 166 MHz or less
- Platform is complete with codec libraries
  - H.264 decode, baseline + main
  - VC-1 decode, baseline + main
  - MPEG-2, MP@ML
  - MPEG-4, SP to ASP L4
  - JPEG, M-JPEG encode
  - AAC, MP3 12-channel encode
- Linux or standalone operation
- ARCAngel™ 4 development platform
  - Virtex-4™ FPGA runs codecs real-time

## ARC™ Media Subsystem Device<sup>1</sup>

ARC 750D core with 32K I-cache,  
32K D-cache

- + SIMD extensions
- + Entropy decode extensions
- + DMA
- + 32K SIMD data memory
- + 10K SIMD code memory



9mm<sup>2</sup>

## Other multicore solutions<sup>2</sup>...



42mm<sup>2</sup>

1. Synthesized in 0.13μm technology for operation at 500 MHz

2. Scaled to 0.13μm for direct comparison, from 81mm<sup>2</sup> at 0.18μm, drawn approximately to scale



- SIMD extensions deliver high performance, yet low power
- Extends ARC™ 700 core family with 104 new instructions
- Decoupled programming model maximizes parallelism
  - Natural and simple extension to single-thread model
- Media-oriented DMA shoulders data transfers
- Supports programmable multistandard audio / video codecs
  - Competitive with dedicated hardware solutions
- Supported by full suite of ARC and GNU software tools
- Integrates with third-party IP
- ARC Media subsystem available Q4
- Symmetric multicore
  - ARC-based implementations now
  - New core products in 2006

## ARC™ Media Subsystem Floorplan

ARC 750D core with 32K I-cache, 32K D-cache

- + SIMD extensions
- + Entropy decode (MPEG-2, MPEG-4, H.264, VC-1)
- + DMA
- + 32K SIMD data memory
- + 10K SIMD code memory

Area: 9mm<sup>2</sup> in 0.13µm TSMC, Virage libraries

Speed: 500 MHz (worst-case, LVLK-OD process)

